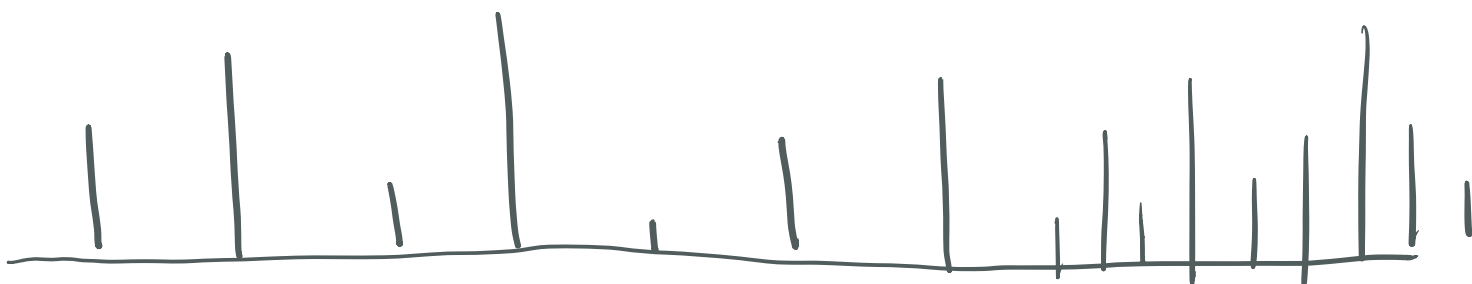


Space efficient quantile selection

Input: stream $a_1, \dots, a_n \in U$ where U has order $<$

(5) (9) (3) (11) (1) (4) (6) 000



e.g. numerical data

names w/ alphabetic order

grades

allowed multiple passes

Goal: return the median w/ minimum:

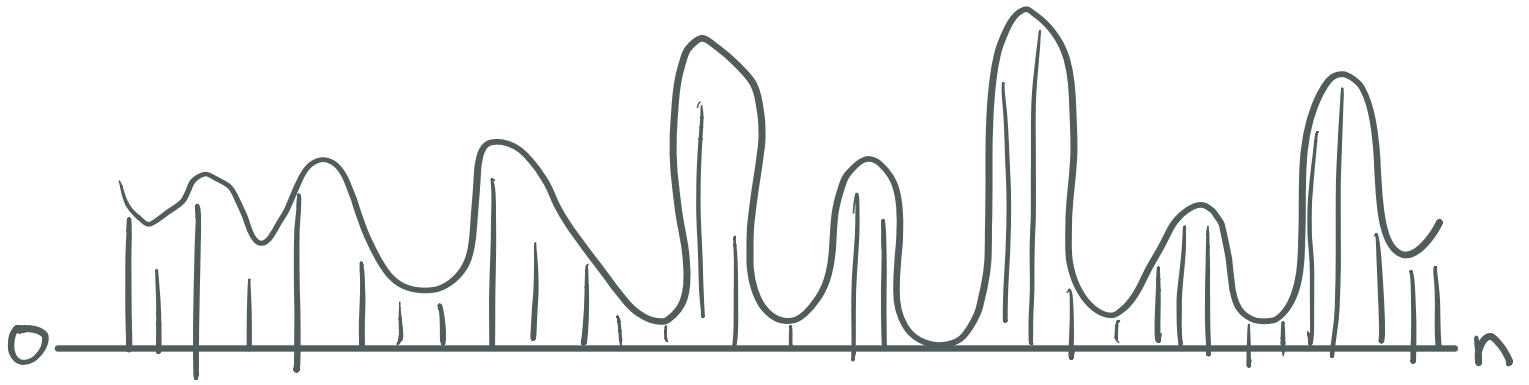
(a) # passes (b) space

more generally: "quantile queries"

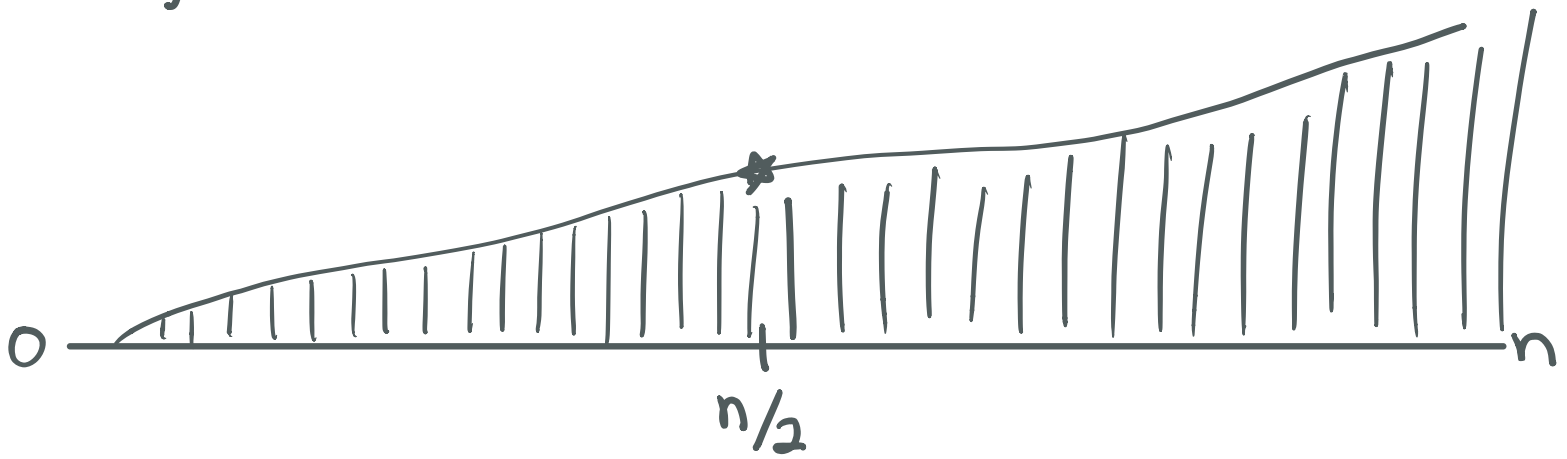
select rank k element

(k th largest)

1 pass: Input



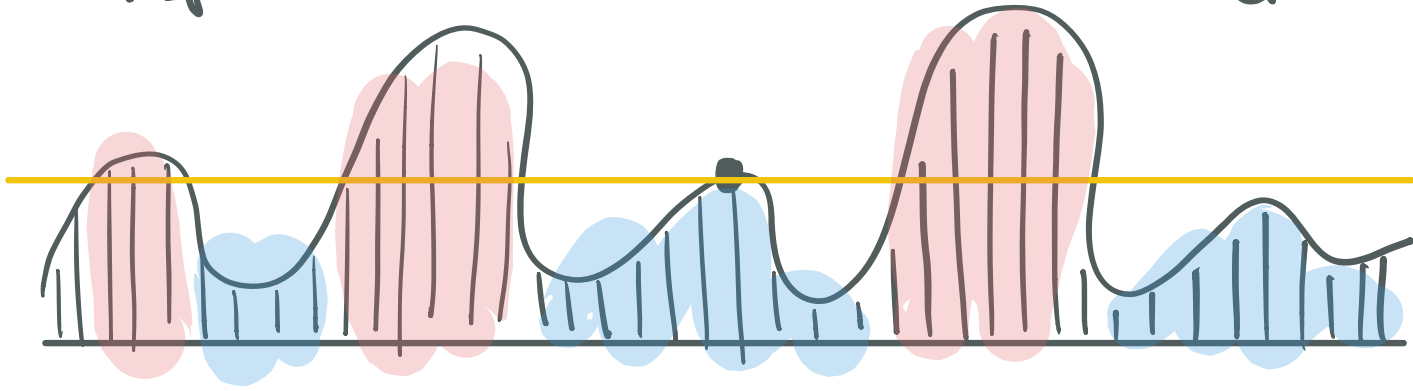
(sort)



sort and select

$O(1)$ space

Quickselect



count # above, below pivot
recurse on appropriate half

<u>Passes</u>	<u>Space</u>	
1	$O(n)$	sort and select
$O(\log n)$	$O(1)$	quickselect (random pivot)
p	$\tilde{O}(n^{1/p})$	Munro, Paterson [1980]



this lecture

Approximations

given rank $k \in [n]$ and param $\epsilon > 0$,
return element w/ rank $k \pm \epsilon n$

Sampling:

for median:

sample $l = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ elements

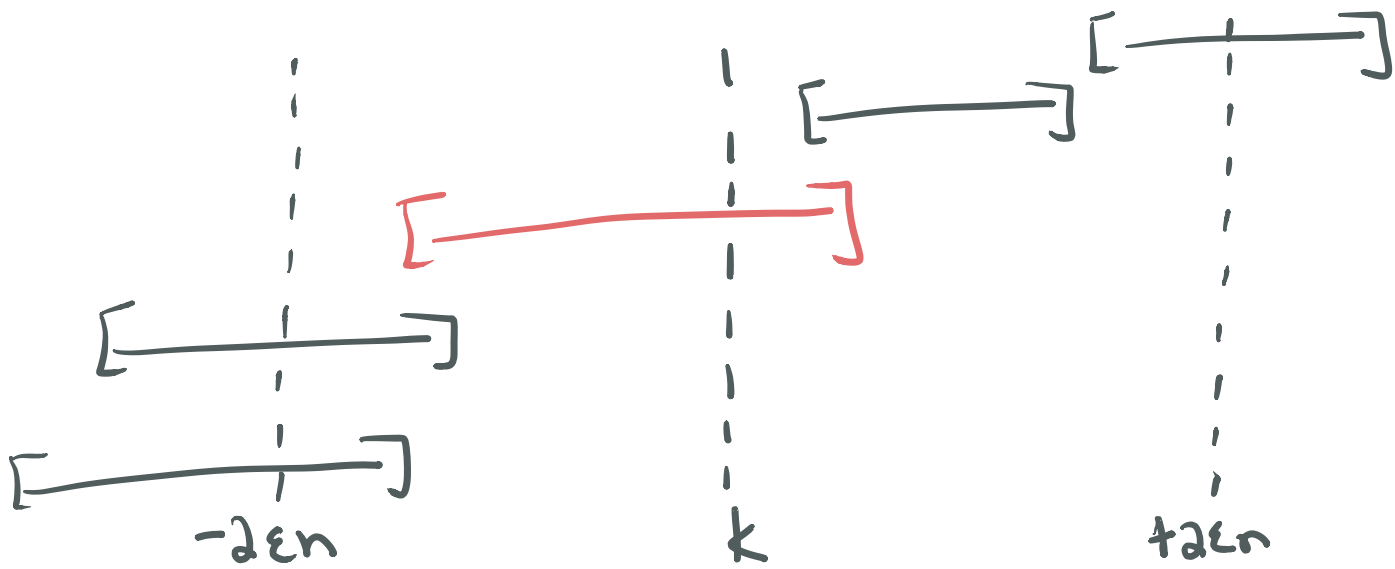
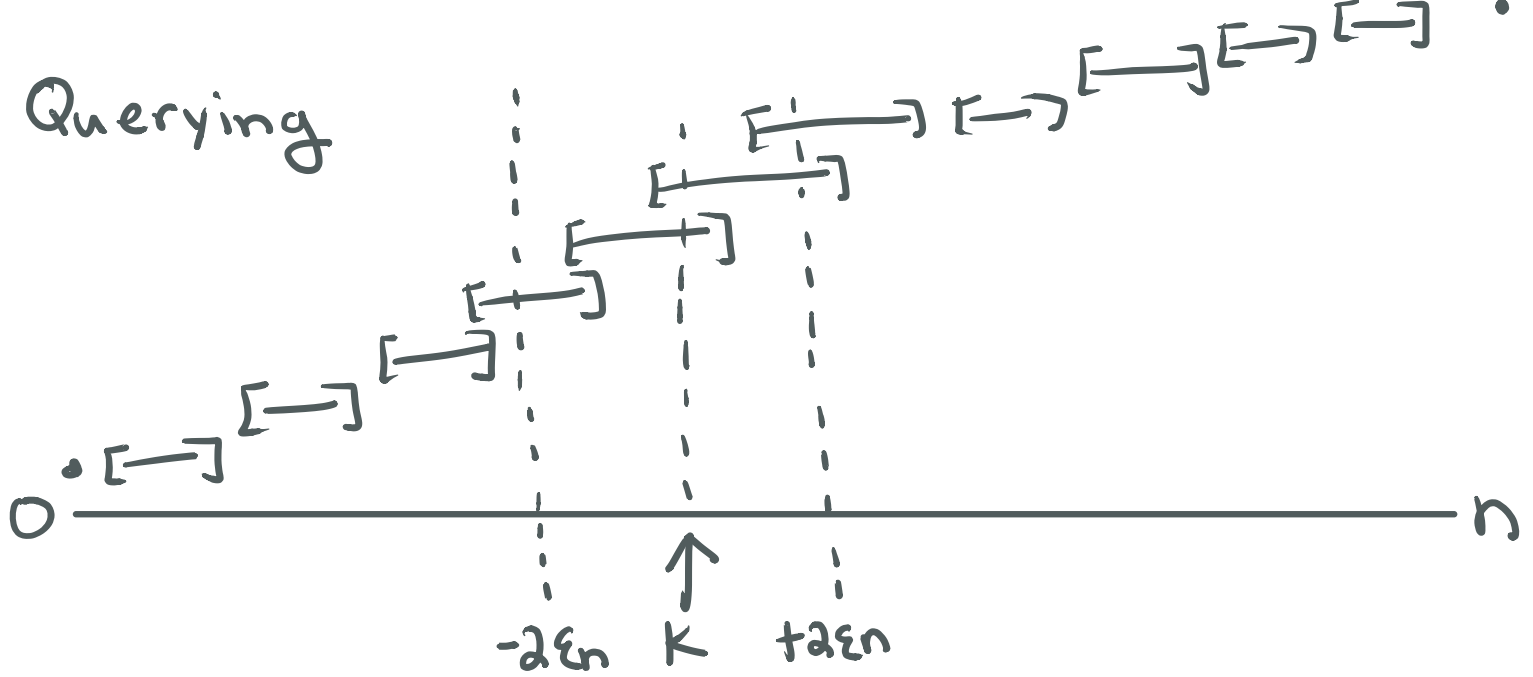
return median of sample

for rank $k = \alpha n$

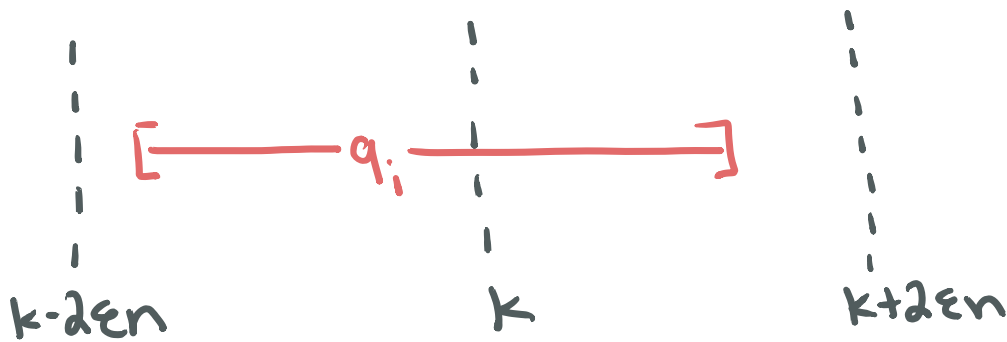
return rank αl element of sample

Deterministic?

Querying

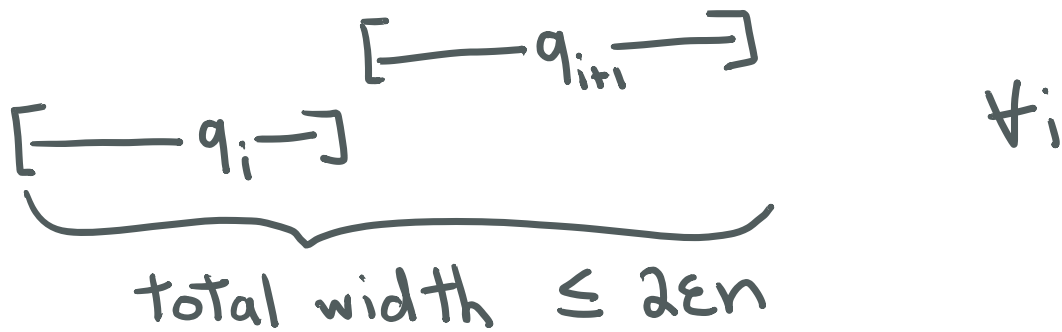


If $I(q_i) \subseteq [k-2\epsilon n, k+2\epsilon n]$ for some q_i ,
then return q_i .



how to ensure such q_i exists $\forall k$?

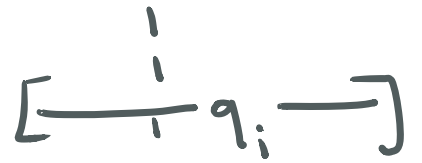
lemma



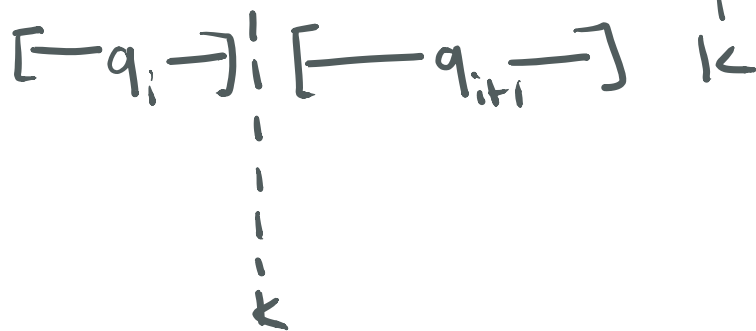
Then every query k contains an interval $I(q_i)$.

Proof two cases:

$k \in I(q_i)$ for some q_i



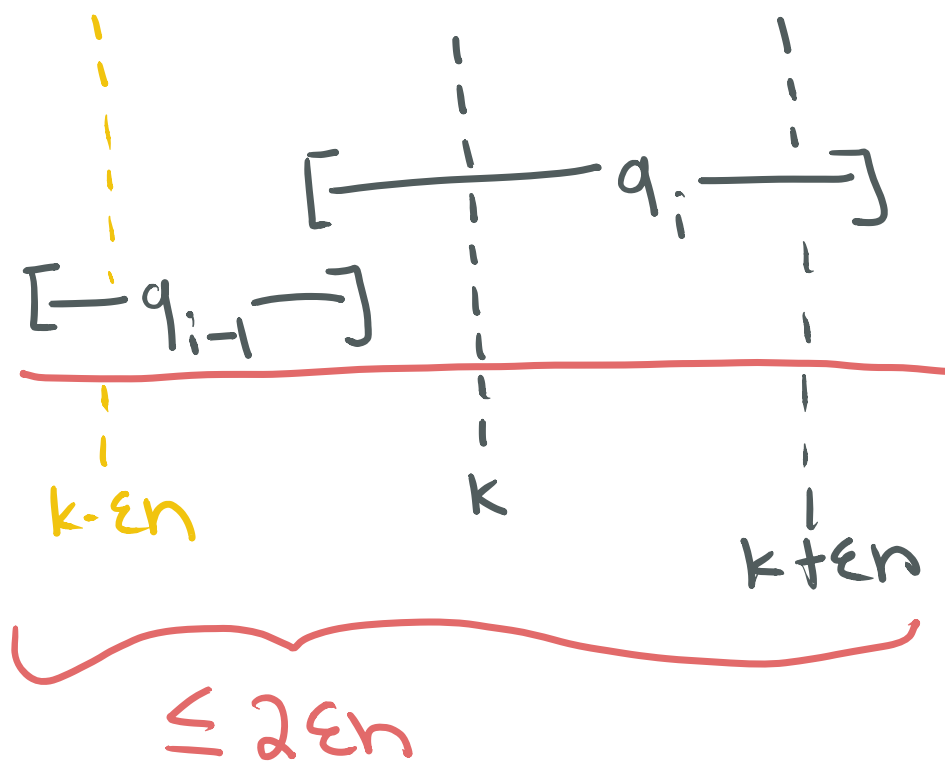
$k \notin I(q_i) \forall q_i$



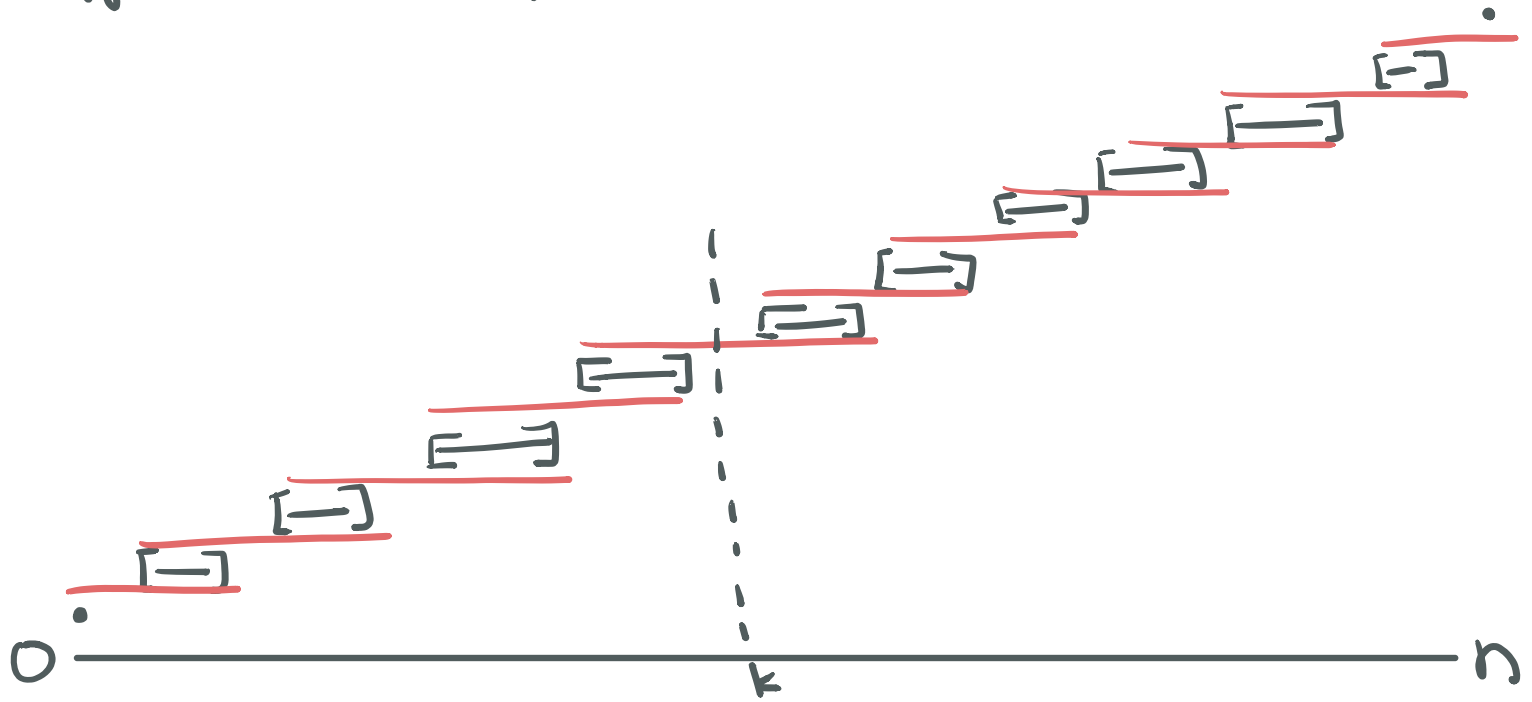
Suppose $k \in I(q_i)$ for some i

if $I(q_i) \subseteq [k-2\varepsilon, k+2\varepsilon]$ then done

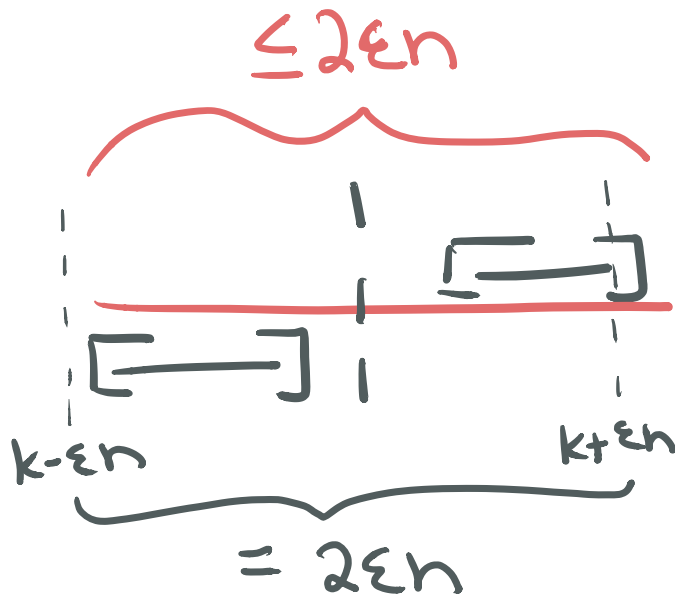
else look at q_{i-1}



suppose $k \notin I(q_i) \forall i$



the "combined intervals" cover $[n]$
pick I covering k .



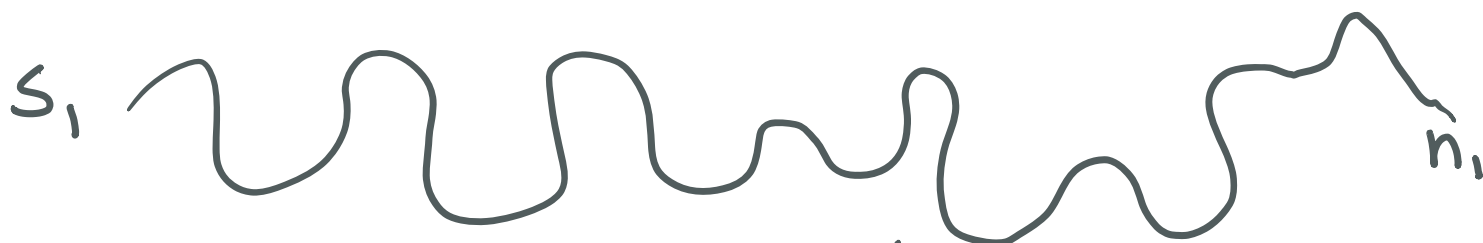
one of the intervals must lie inside

Key invariant: any two consecutive intervals have width $\leq 2\epsilon n$.

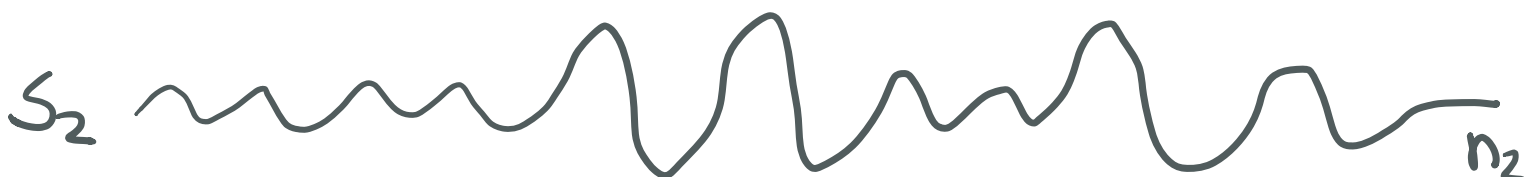
" ϵ -APX quantile summary"

Merging given two ϵ -APX quantile

summaries over 2 streams, want ϵ -APX
summary over combined stream

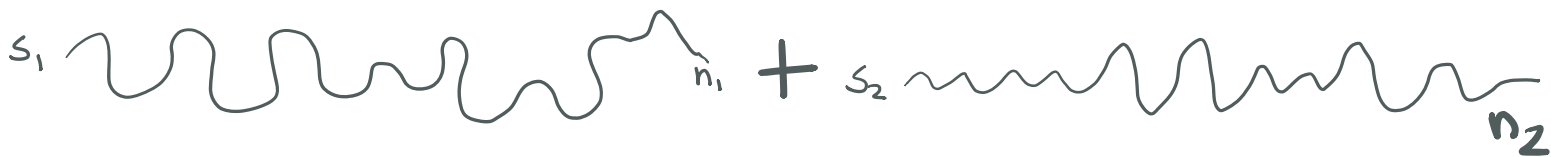


$$\Rightarrow Q' = \left[\begin{array}{c} a'_1 \\ \left[\begin{array}{c} a'_2 \\ \left[\begin{array}{c} a'_3 \\ \left[\begin{array}{c} a'_4 \\ \left[\end{array} \right] \end{array} \right] \end{array} \right] \end{array} \right] \end{array} \right]$$



$$\Rightarrow Q'' = \left[\begin{array}{c} \left[\begin{array}{c} \left[\begin{array}{c} \left[\end{array} \right] \end{array} \right] \end{array} \right] \end{array} \right]$$

want to combine Q', Q'' to get summary of



$$"Q' \cup Q''" = \{ a''_1, \dots, a''_c, I''(a''_1), \dots, I''(a''_c) \}$$

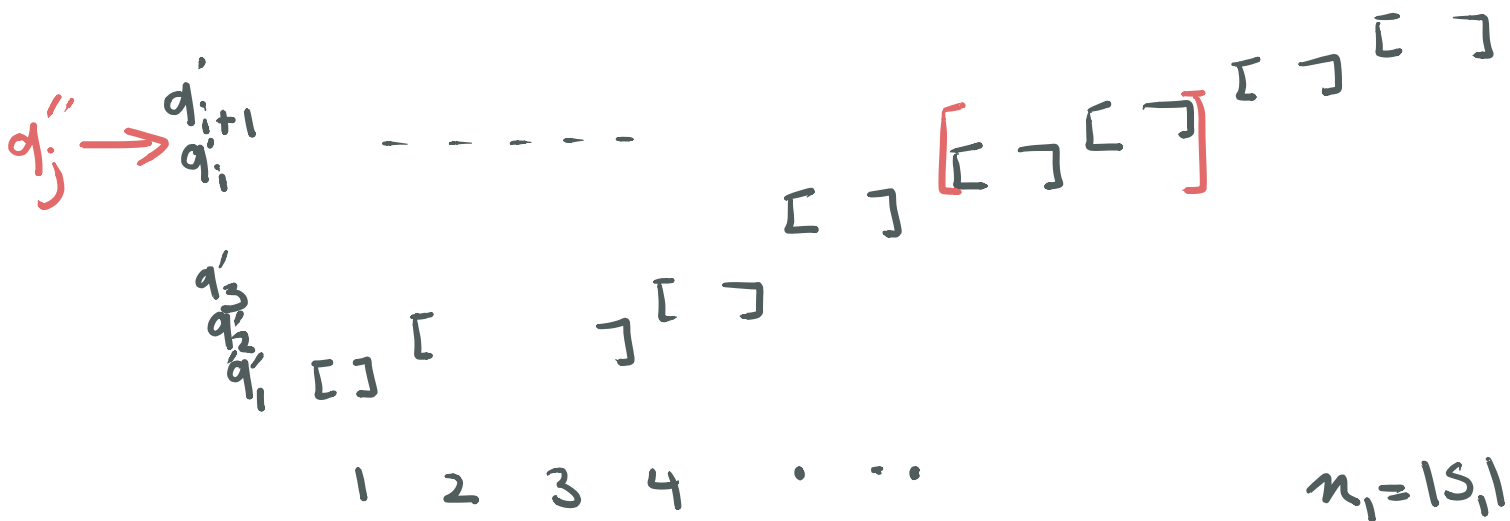
denote:

$$Q' = \{q'_1, \dots, q'_\ell, I'(q'_1), \dots, I'(q'_\ell)\}$$

$$Q'' = \{q''_1, \dots, q''_m, I''(q''_1), \dots, I''(q''_m)\}$$

let $q''_j \in Q''$. $I''(q''_j)$ bounds rank q''_j w.r.t S_2

goal: bound rank q''_j w.r.t $S_1 + S_2$.



$$\begin{aligned} \text{rank of } q''_j \text{ in } S_1 &\geq \min I'(q_i) \\ &\leq \max I'(q_{i+1}) \end{aligned}$$

$$\text{so } \min I'(q_i) + \min I''(q''_j)$$

$$\leq \text{rank}(q''_j | S_1 + S_2)$$

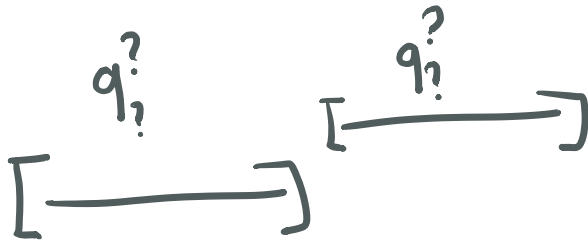
$$\leq \max I''(q''_j) + \max I'(q_i)$$

$$\text{set } I'''(q''_j) = \left[\min I'(q_i) + \min I''(q''_j), \max I''(q''_j) + \max I'(q_i) \right]$$

$$Q''' = \{q'_1, \dots, q'_e, q''_1, \dots, q''_m, \text{ w/ intervals } I''_1\}$$

to show Q''' is ϵ -APX, need to show
"2 ϵ n width" property.

Take two consecutive intervals in Q''' .

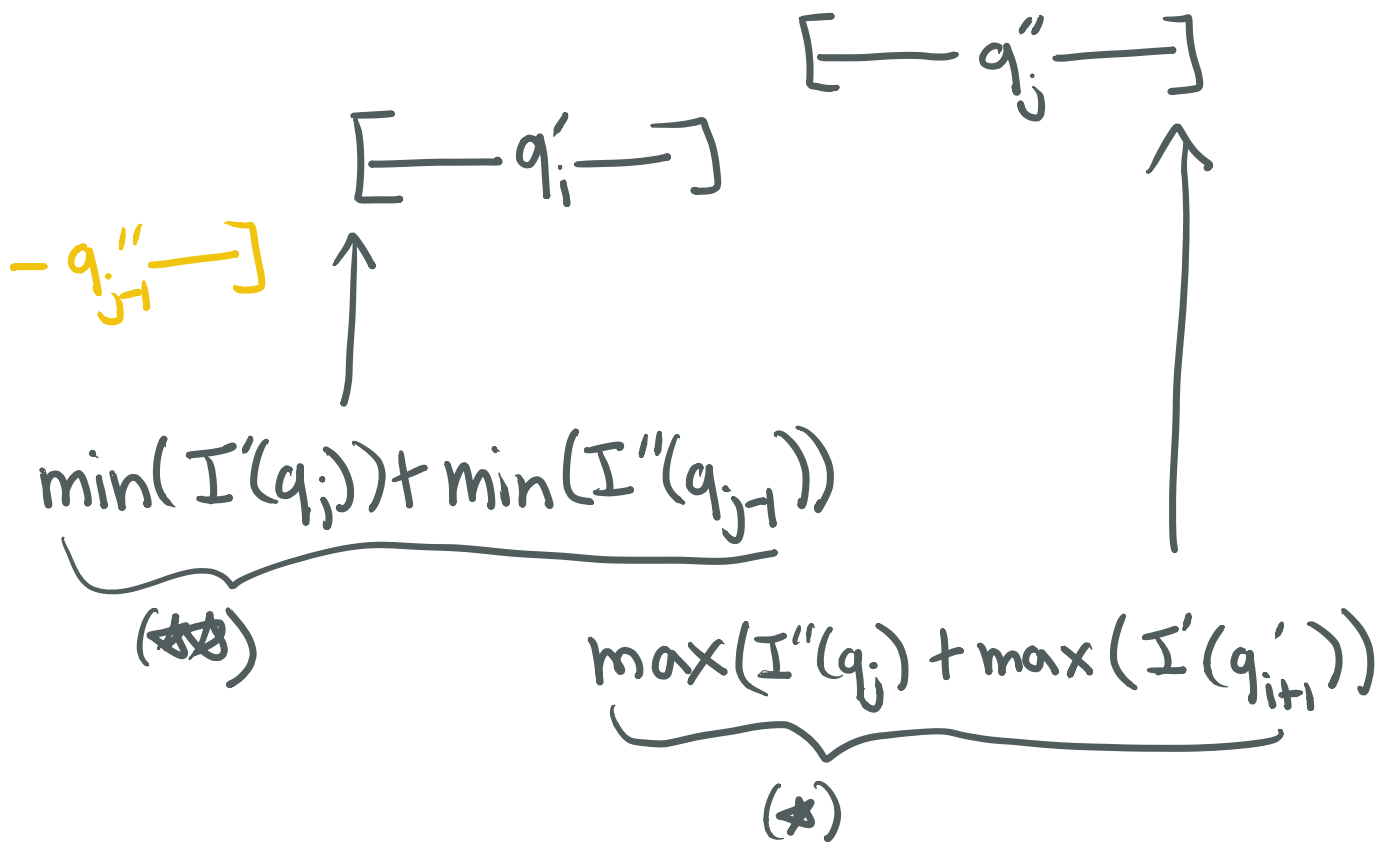


Two cases:

- ① elements from diff sets
- ② elements from same sets

① diff sets:

$[\text{---} q'_{i+1} \text{---}]$

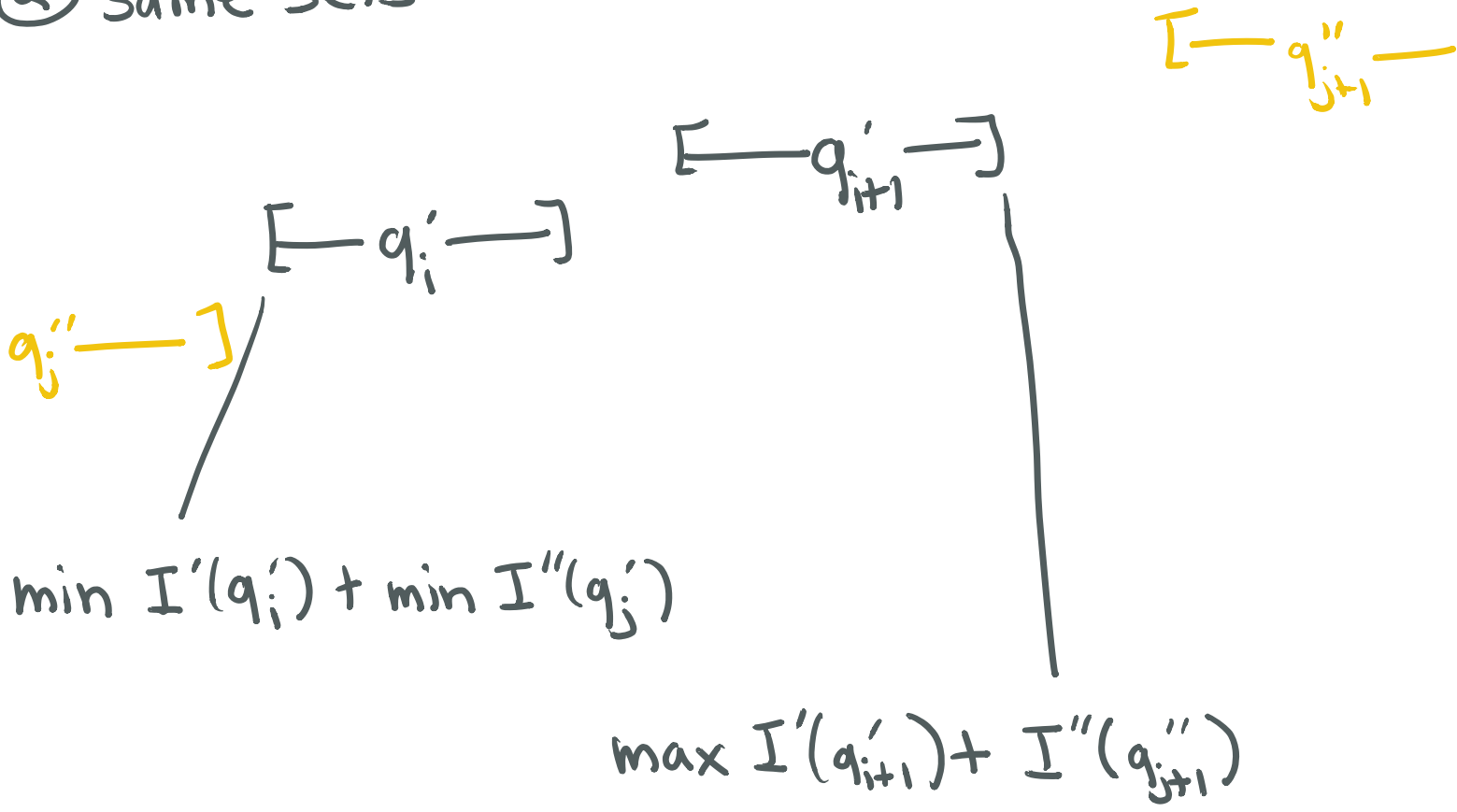


$$\begin{aligned} (*) - (**) &= \max(I'(q_{i+1}) - \min I'(q_j)) \leq 2\varepsilon n_1 \\ &+ \max(I''(q_j)) - \min I''(q_{j-1}) \leq 2\varepsilon n_2 \end{aligned}$$

$$\leq 2\varepsilon(n_1 + n_2)$$



② same sets



$$\frac{\max I'(q_{i+1}') + I''(q_{j+1}'') - \min I'(q_i') + \min I''(q_j'')}{\epsilon n_1 + \epsilon n_2 = \epsilon(n_1 + n_2)}$$



This shows that merging 2 ϵ -APX QS's gives ϵ -APX QS of combined streams

Lemma

given ϵ -approximate quantile summaries

Q_1, \dots, Q_h for S_1, \dots, S_n , we can

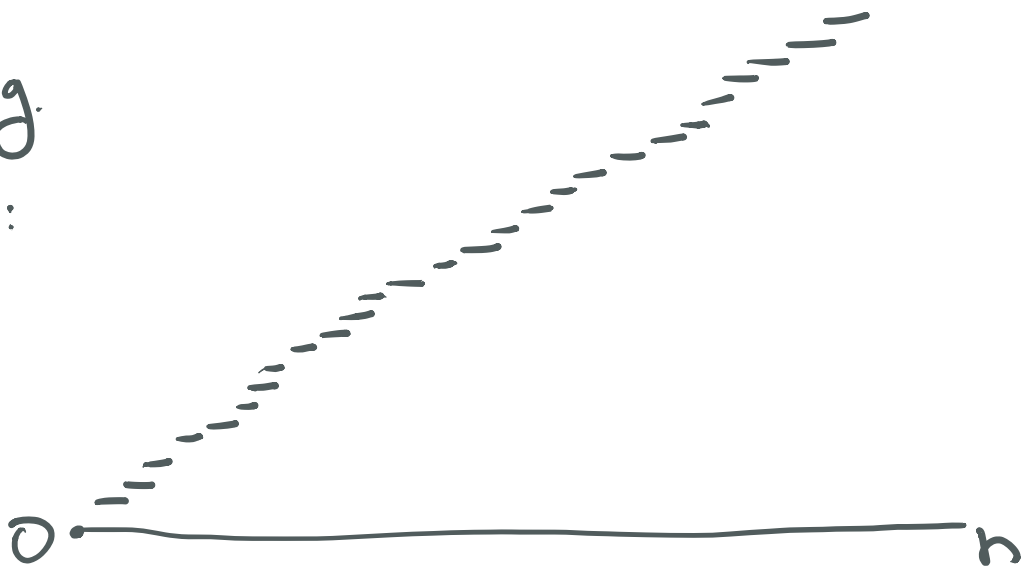
combine to make ϵ -approximate summary

$Q_1 \cup \dots \cup Q_h$ of $S_1 \cup \dots \cup S_n$.

Size?

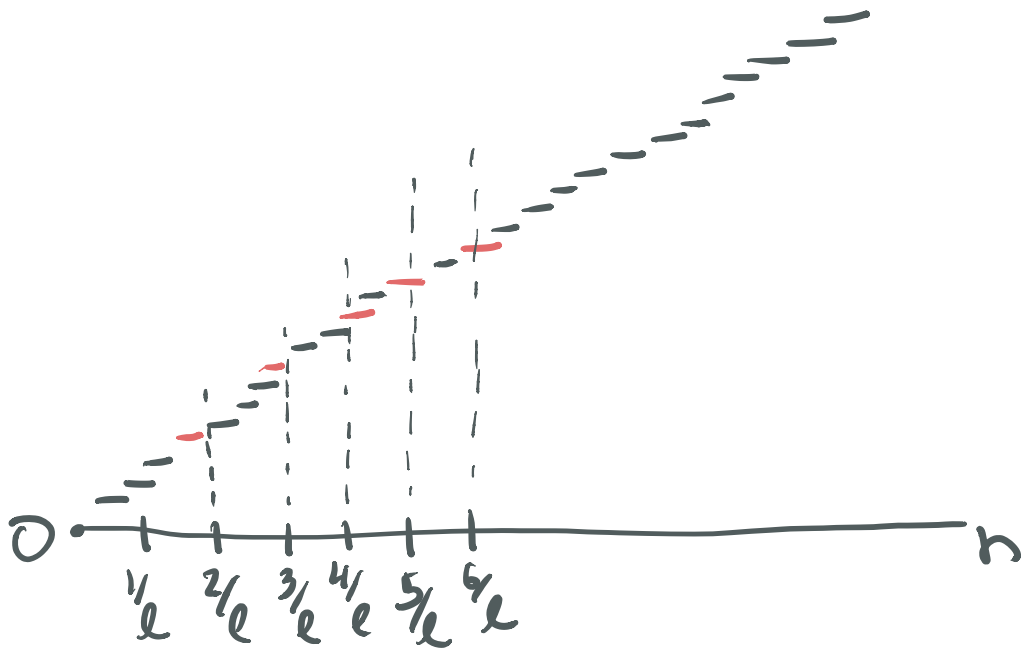
Pruning.

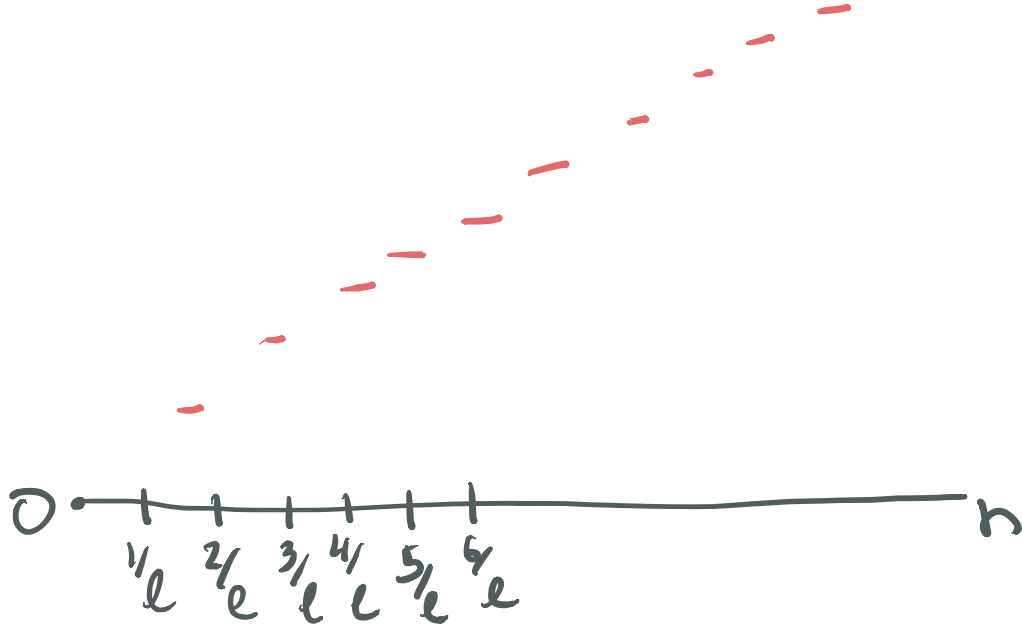
Input:



ϵ -approximate quantile summary w/ too many points

Goal: sparser summary that's still very good



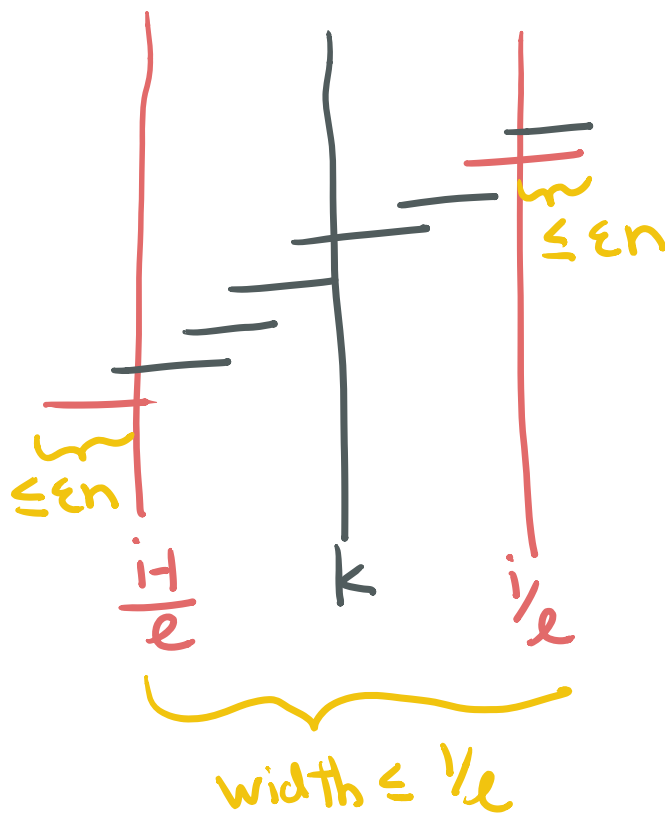


claim: resulting quantile is $(\epsilon + \frac{1}{2\epsilon})$ -APX

Proof

suppose we query a particular rank k

look at original summary



total width $\leq 2\epsilon n + \frac{1}{e}$ i.e. $(\epsilon + \frac{1}{2\epsilon})$ -approx.

Recap:

we can combine ϵ -APX quantile summaries to get ϵ -APX quantile summary of whole thing

sparsify ϵ -APX quantile summary to $(\epsilon + \frac{1}{2k})$ -APX quantile summary w/ k points

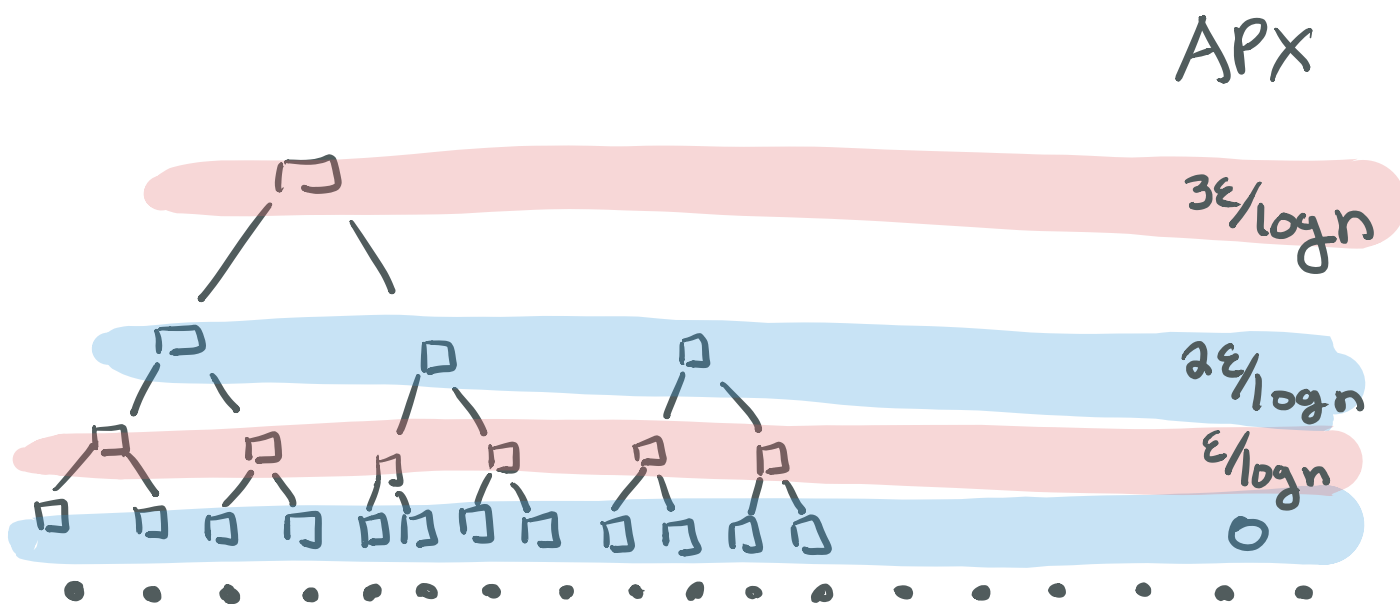
Remains to address:

how to make one at all??

what if $n=1$?

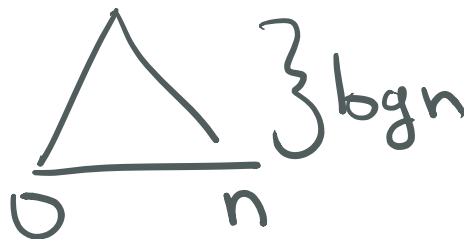
make a summary of just that point.

I claim that's all we need!



take $k = \frac{\ln n}{\epsilon}$

at the root,

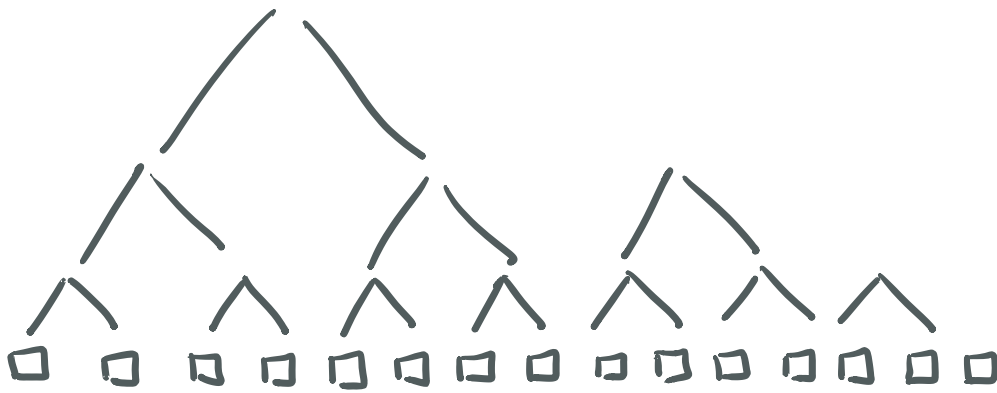


ϵ -approximate quantiles

Space?

each summary takes $O(\log(n)/\epsilon)$ space.

How many summaries?



only keep "root summaries"

how many "root summaries"?

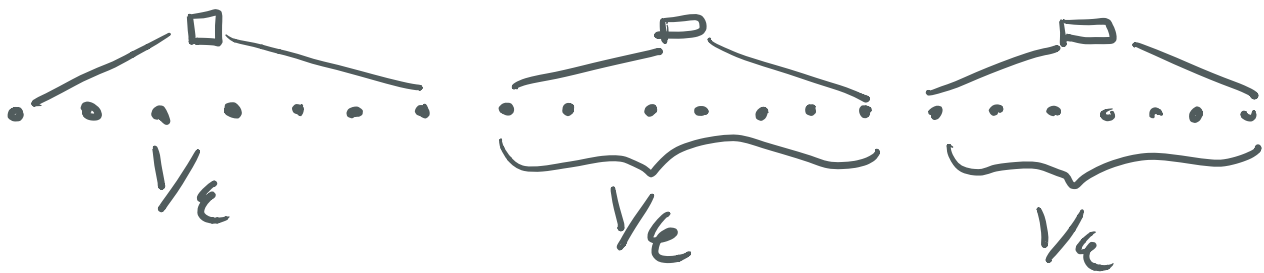
One summary per height, so $\log n$ total summaries at any time

Theorem

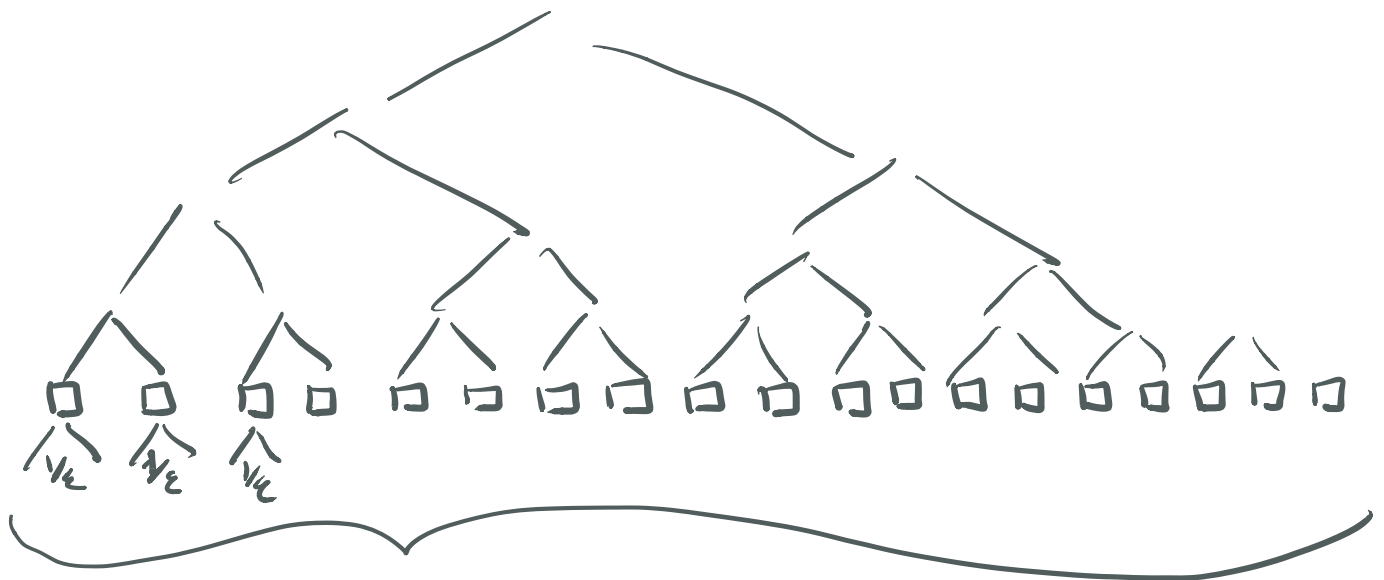
- 1-pass
- $O(\log^2(n)/\epsilon)$ space
- deterministic
- ϵ -APX quantile over stream

idea: mergability + dyadic intervals trick

slightly better:



have first level contain $1/\epsilon$ points



ϵn summaries at "leaf level"

$\Rightarrow \log(\epsilon n)$ height

Theorem ††

- 1-pass
- $O(\log^2(\epsilon n)/\epsilon)$ space
- deterministic
- ϵ -APX quantile over stream

Even better?

Khanna-Greenwald:

$\frac{1}{\epsilon} \log(\epsilon n)$ space

- more sophisticated quantile summary, merging
- interval trick

Finding the median (and other ranks)
in p passes

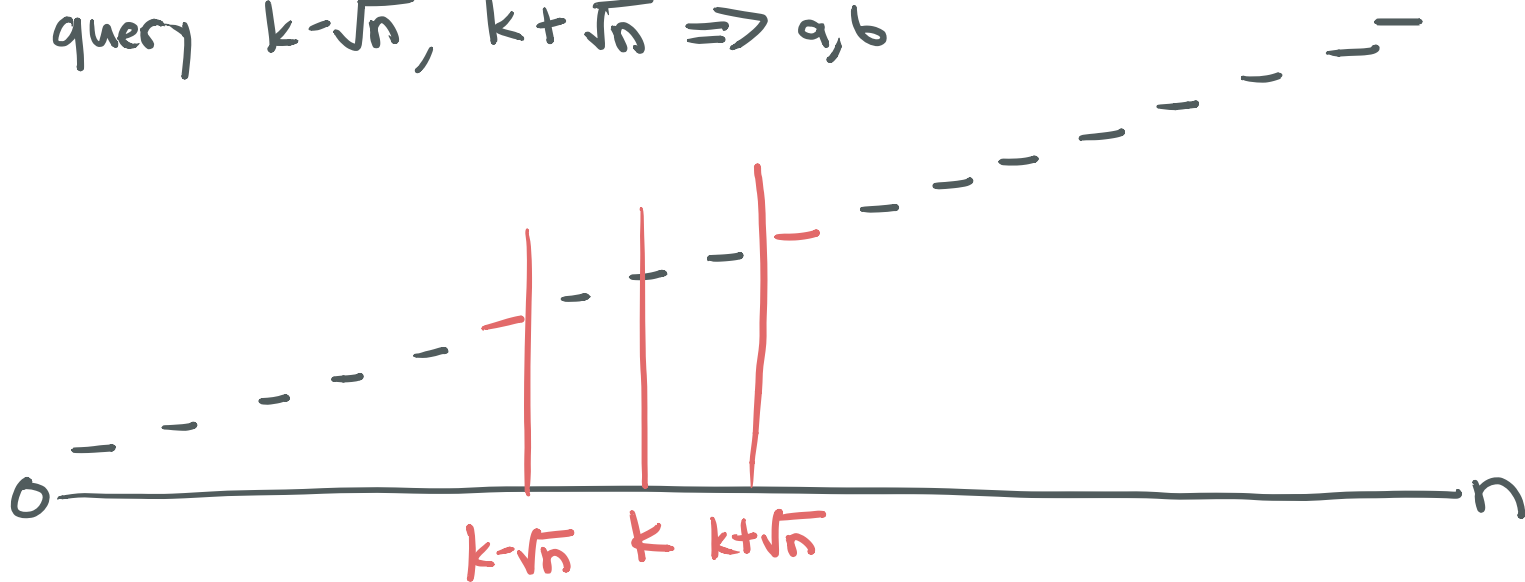
Fix $p=2$ for simplicity.

goal: $O(\sqrt{n} \text{ polylog}(n))$ space

suppose we are querying rank k .

1st pass: build ϵ -APX quantile summary
for $\epsilon = 1/\sqrt{n}$ ($\sqrt{n} \log(n)$ space w/ GK)

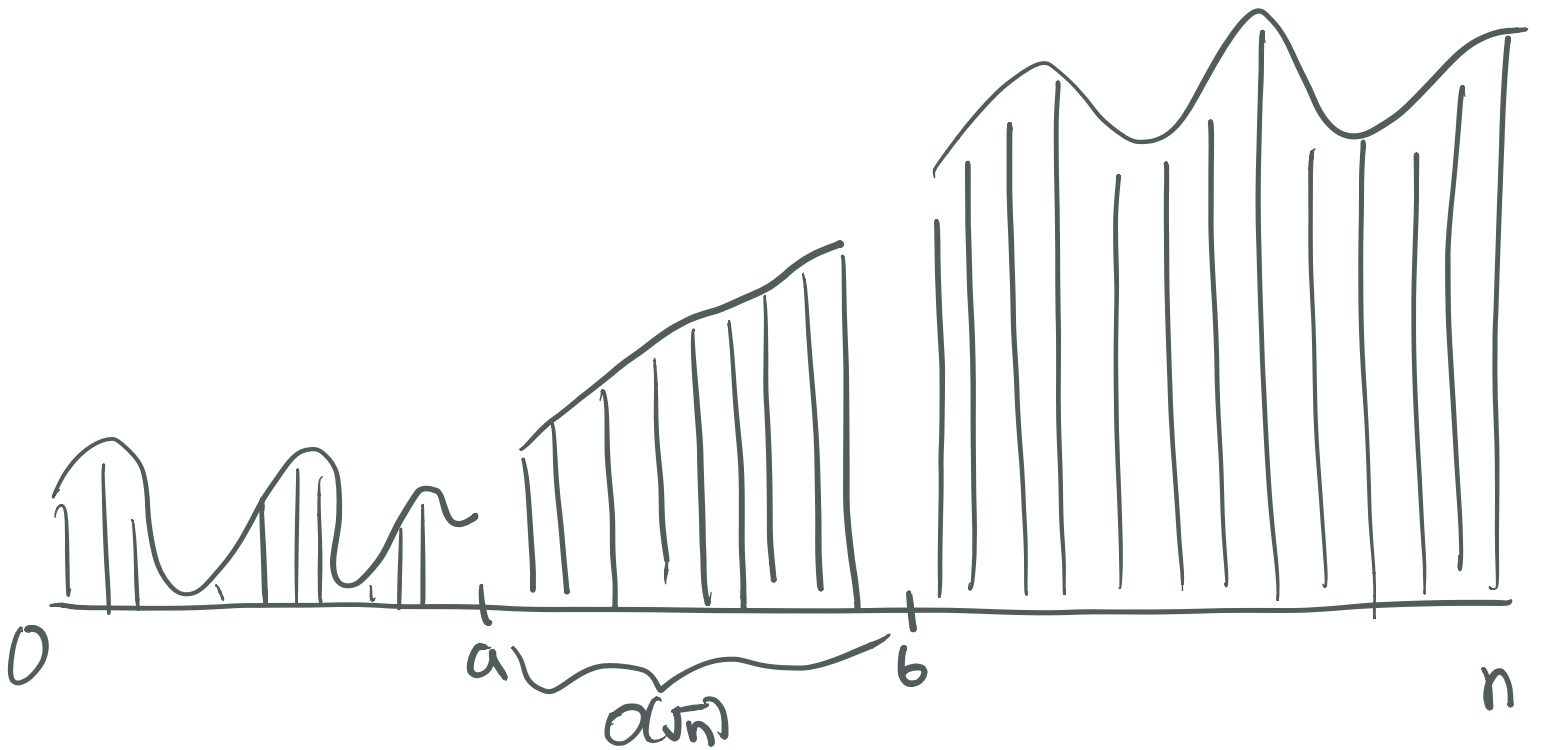
query $k - \sqrt{n}$, $k + \sqrt{n} \Rightarrow a, b$



$$k - 2\sqrt{n} \leq \text{rank}(a) \leq k \leq \text{rank}(b) \leq k + 2\sqrt{n}$$

2nd pass:

- take all elements between a and b



- compute $\text{rank}(a)$, (by counting)
- sort taken elements
- return $(k - \text{rank}(a))$ th taken element

for general p :

make $\frac{1}{n^{1/p}}$ -APX quantile summaries
and filter.

After p passes, down to $n^{1/p}$ elements
sort and select.